



Online Audio-Visual Autoregressive Speaker Extraction

Zexu Pan¹, Wupeng Wang², Shengkui Zhao¹, Chong Zhang¹, Kun Zhou¹, Yukun Ma¹, Bin Ma¹

¹Tongyi Lab, Alibaba Group, Singapore

²National University of Singapore, Singapore

zexu.pan@alibaba-inc.com

Abstract

This paper proposes a novel online audio-visual speaker extraction model. In the streaming regime, most studies optimize the audio network only, leaving the visual frontend less explored. We first propose a lightweight visual frontend based on depth-wise separable convolution. Then, we propose a lightweight autoregressive acoustic encoder to serve as the second cue, to actively explore the information in the separated speech signal from past steps. Scenario-wise, for the first time, we study how the algorithm performs when there is a change in focus of attention, i.e., the target speaker. Experimental results on LRS3 datasets show that our visual frontend performs comparably to the previous state-of-the-art on both SkiM and ConvTasNet audio backbones with only 0.1 million network parameters and 2.1 MACs per second of processing. The autoregressive acoustic encoder provides an additional 0.9 dB gain in terms of SI-SNRi, and its momentum is robust against the change in attention.

Index Terms: Cocktail party problem, real-time, speaker extraction, visual frontend, autoregressive

1. Introduction

Real-world speech signals are often mixed with interfering speech and noise signals. Human brains excel at focusing on an interested speech signal, i.e., target speech, while filtering out the rest, known as auditory attention [1]. Equipping machines with such auditory selective attention is crucial for speech applications such as automatic speech recognition [2].

Speech separation algorithms separate a multi-talker speech signal into individual streams by speakers [3–7]. However, the separated speech signals are not associated with specific speaker identities. In contrast, speaker extraction research exhibits the attention to disentangle only a target speaker’s speech signal from the multi-talker speech signal, with the attention driven by an auxiliary conditioning signal. The auxiliary signal could be a pre-recorded speech signal for the network to attend to a speech that has a similar voice signature [8–10], or a face recording for the network to attend to the synchronized speech [11–15], or even the brain signal for the network to attend to the speech that the brain is trying to listen [16–18].

This research focuses on audio-visual speaker extraction (AVSE) in online scenarios, addressing the need for streaming on-device AVSE algorithms in applications like video conferencing and human-robot interactions, such as in-car or service robot communication. Due to device constraints, on-device models usually favor fewer network parameters and multiply-accumulate (MAC) operations, to accelerate the processing, which also enhances the real-time factor. While most existing research emphasizes optimizing the audio component of

the network [19–23], limited studies have focused on improving the visual counterpart. The visual encoder is still predominantly based on ResNet18 [23, 24], with recent explorations into ShuffleNet-based networks [25]. In this work, inspired by the effectiveness of the BlazeFace face detection network [26] and depth-wise separable convolution [27], we propose a lightweight visual encoder that reduces feature dimension size while incorporating deeper depth-wise separable convolutional layers to encode lip motion more efficiently.

We also aim to enhance the audio components of online AVSE. During real-time inference, while the network processes newly acquired frames of a mixture signal, it also has access to the separated speech signals from previous frames. This availability makes autoregressive networks well-suited to leveraging this past information. In the literature, the streaming NeuroHeed model, which conditions on brain signals, encodes the separated signals into a single speaker embedding vector using a speaker encoder and repeatedly fuses it with every speech frame in its sliding window-based decoding process [17]. The PARIS speech separation research also explored incorporating past separated speech frames into the mixture signal inputs to utilize historical information in frame-level decoding [28]. In this work, we adopt frame-level decoding like PARIS, which is better suited for online decoding with improved real-time performance, differently, we propose encoding the past extracted signals with a lightweight acoustic encoder, to generate distinct acoustic embeddings at the frame level for the speaker extractor, which serves as a secondary conditioning signal alongside the visual cue.

Another common scenario in real-world deployments of streaming AVSE models is the changing focus between active speakers, such as in multi-party meetings and conversations. The network must adjust its attention dynamically to extract the correct speech signal. Previous works assume the identity of the target speaker is consistent. In this work, we explore the changing target scenario for the first time, demonstrating that visual-conditioned speaker extraction models exhibit robust performance when faced with such changes in attention, which is crucial for real-world applications. Furthermore, our proposed acoustic encoder, which encodes past information into frame-level conditioning embeddings, shows robustness in adapting to changes in attention. This contrasts with the NeuroHeed model, where aggregating all information into a single speaker embedding may create momentum that hinders the network’s ability to adjust focus quickly.

The contributions of this work for online AVSE studies are threefold: 1) We propose a lightweight visual encoder for AVSE, which, with only 0.1 million parameters and 2.1 MACs per second of processing, performs comparably to the previous state-of-the-art visual encoders on both SkiM and ConvTas-

Net audio backbones. 2) We introduce an autoregressive acoustic encoder for streaming AVSE, which provides an additional 0.9 dB gain in SI-SNR_i on the Lip Reading Series 3 (LRS3) dataset [29]. 3) We conduct the first study on the switching attention scenario, demonstrating that our proposed model is robust in such contexts.

2. Proposed network

Let $x(\tau)$ be a multi-talker mixture speech signal¹, consisting of the target speech signal $s(\tau)$ and interference speech signal $b(\tau)$, the AVSE network $f(\cdot)$ estimates the target speech signal $\hat{s}(\tau)$ to approximate $s(\tau)$, conditioned on the visual recording of the target speaker $v(t)$:

$$\hat{s}(\tau) = f(x(\tau), v(t)) \quad (1)$$

In this work, we additionally study the changing target scenario, meaning that the identity of $v(t)$ and $s(\tau)$ may change in a processing clip. For simplicity, we maintain the same representation of $v(t)$ in this changing target scenario.

2.1. Architecture

Our proposed online audio-visual autoregressive speaker extraction network is illustrated in Fig. 1. It comprises two identical non-shared-weight speech encoders, a speaker extractor, a speech decoder, a novel visual encoder, and a novel autoregressive acoustic encoder.

For speech encoder and decoder, we follow the masking-based time-domain approach [4, 7, 21]. The speech encoder includes a 1-dimensional (1D) convolution *Conv1D* and a rectified linear activation. The speech decoder consists of a linear layer and an overlap-and-add operation. The channel, kernel, and stride sizes are set to 128, 16, and 8, respectively.

For the speaker extractor, we adopt the online skipping memory long short-term memory network (SkIM), which is known for its low latency and excellent performance for online speech separation [21]. We set the number of output streams to one to suit the speaker extraction task. We set the long short-term memory (LSTM) unit size to 384, the number of layers to 3, and the non-overlapping segment size to 50.

Our proposed lightweight visual encoder, named BlazeNet64, is shown in the top section of Fig. 1, it encodes the lip image sequence $v(t)$ to visual cue, which is used as an attention attractor to extract the target speech. The architecture is motivated by the effectiveness of the BlazeFace face detection network [26] and depth-wise separable convolution [27]. $v(t)$ is first processed by a causal 3-dimensional convolution *Conv3D*. The rest of the layers in the visual encoder are depth-wise separable 2-dimensional convolutions applied to each image independently. Compared to the widely used *ResNet18*-based visual encoder, which was originally designed for lip reading [30], our visual encoder is narrower but deeper, and much more light-weighted. We do not pre-train this visual encoder, so we do not use adapter layers after the visual encoder as in [13, 31]. To align with the audio embeddings, the visual embeddings $V(t)$ are temporally repeated to account for frame rate differences before concatenation with the audio embeddings at the beginning of the speaker extractor.

Our proposed lightweight autoregressive acoustic encoder leverages extracted clean speech samples from past frames to

¹Variables with τ denote speech signals in the time domain, while variables with t denote frame-based embeddings.

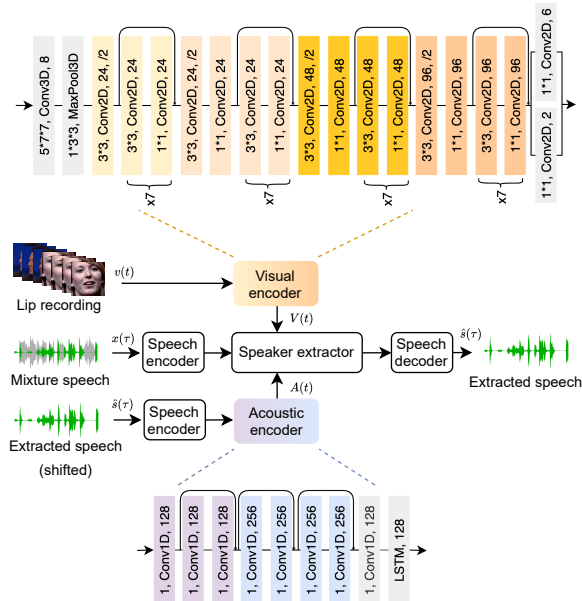


Fig. 1: Our proposed online audio-visual autoregressive speaker extraction network. The middle section introduces the network components and data flow, the top section showcases our proposed lightweight visual encoder, and the bottom section depicts our proposed lightweight autoregressive acoustic encoder. Each colored layer’s parameters are specified in terms of kernel size, layer type, output channel size, and stride.

maintain temporal attention momentum during online processing, shown in the bottom section of Fig. 1. The extracted speech samples from the past frames will first pass through a speech encoder, and then through a series of 1D-convolutional layers before being summarized by a single LSTM layer. The output $A(t)$ is also concatenated with the audio embeddings at the start of the speaker extractor.

2.2. Training strategy and loss function

The network is non-autoregressive if the acoustic encoder is dropped from Fig. 1, in this case, the negative scale-invariant signal-to-noise ratio (SI-SNR) [32] is adopted as the loss function for training:

$$\mathcal{L}_{\text{SI-SNR}}(s, \hat{s}) = -20 \log_{10} \frac{|\langle \hat{s}, s \rangle|}{|\hat{s}| |s|} \frac{|s|}{|\hat{s} - \frac{\langle \hat{s}, s \rangle}{|s|^2} s|}. \quad (2)$$

When our proposed acoustic encoder is integrated into the overall network, the network operates in an autoregressive mode. In this case, we adopt the Pseudo-AutoRegressive Siamese Training (PARIS) strategy [28], which involves two forward passes for each batch during training. The first pass forward the network without the acoustic encoder (with $A(t)$ being set to zero embeddings as a placeholder). The second pass performs a full forward pass with the input of the acoustic encoder being the output of the first pass which is used as the “pseudo past extracted speech”.

The PARIS paper suggested using the scale-sensitive signal-to-noise ratio (SNR) rather than the widely used SI-SNR as the loss function for training, to normalize network output for feedback in step-by-step streaming settings. However, networks trained with the SNR loss function typically do not perform as well as those trained with the SI-SNR loss function.

Therefore, in this work, we pair SI-SNR with a scale-variant frequency domain loss [33], specifically the frequency-domain multi-resolution delta spectrum loss $\mathcal{L}_{\text{freq}}$ [33], as the loss function. The loss functions for the first pass \mathcal{L}_1 and the second pass \mathcal{L}_2 are defined as:

$$\mathcal{L}_1 = \mathcal{L}_{\text{SI-SNR}}(s, \hat{s}^1) + 0.25 * \mathcal{L}_{\text{freq}}(s, \hat{s}^1) \quad (3)$$

$$\mathcal{L}_2 = \mathcal{L}_{\text{SI-SNR}}(s, \hat{s}^2) + 0.75 * \mathcal{L}_{\text{freq}}(s, \hat{s}^2) \quad (4)$$

where \hat{s}^1 and \hat{s}^2 are the network outputs of the first and second pass respectively.

3. Experimental setup

3.1. Dataset

We mainly use the Lip Reading Sentences 3 (LRS3) dataset to validate our proposed method in this work [29], which is widely used in many AVSE studies [34–36]. The speech signal is available at 16kHz, and video is available at 25 frames per second.

We study the following three scenarios in this paper: 1) Two speakers speak simultaneously, with one being the target speaker. 2) Two speakers speak simultaneously, with the target speaker switching from one to the other at a random time. 3) One speaker speaks for a duration, followed by another speaker, both of whom are target speakers. An additional interference speaker overlaps with the target’s speech.

To create the mixture utterances, we first select one utterance as the anchor. The remaining utterances are normalized to match the energy level of the anchor. Each utterance is then assigned a random SNR ranging from 10dB to −10dB relative to the anchor before being combined. We simulate 40,000, 5,000, and 3,000 such two-speaker overlapping mixture utterances for our train, validation, and test sets using the original LRS3 datasets. The speaker and utterances in the test set do not overlap with the train and validation sets to ensure speaker-independent analysis.

3.2. Optimization

We use the PyTorch framework to conduct our experiments². All models trained in this work adopt the same optimization setting. We used distributed data-parallel training with four Tesla 16 GB V100 GPUs, the effective batch size is 16. We train the models for 150 epochs, the Adam [37] optimizer is used with an initial learning rate of 0.001. The learning rate is reduced by half if the best validation loss (BVL) does not improve for 6 consecutive epochs, with training stopping early if the BVL does not improve for 20 consecutive epochs.

4. Results

To evaluate the quality of extracted speech, we use several metrics: the improvement in SI-SNR (SI-SNRi) [32], the improvement in SNR (SNRi) [38], the improvement in Perceptual Evaluation of Speech Quality (PESQi) [39], and the improvement in Short-Term Objective Intelligibility (STOIi) [40]. All improvements are calculated relative to the unprocessed multi-talker speech signals. Higher values indicate better performance. We use SI-SNRi as our main metric when comparing the results, as other metrics show similar trends. For clarity, each differently trained system is assigned a unique system number (Sys).

²The model and training scripts are available at <https://github.com/modelscope/ClearerVoice-Studio>

Table 1: Comparison of visual encoders. Parameters (Param) are reported in million (m), and multiply-accumulate (MAC) operations is in billion (G) per 1 second input, the lower the better.

Visual Encoder	Params (M)	MACs (G)
ResNet18	11.2	12.9
ShuffleNetV2	0.9	7.3
BlazeNet64 (Ours)	0.1	2.1

4.1. Comparison on visual encoders

We first compare our visual encoder with baseline visual encoders in Table 1 and Table 2. We consider two baselines: the ResNet18-based visual encoder and the ShuffleNetV2-based visual encoder, both are state-of-the-art approaches.

In Table 1, we examine the computational efficiency. Our BlazeNet64 visual encoder stands out with only 0.1 million network parameters, which is significantly smaller compared to the 0.9 million and 11.2 million parameters of the ShuffleNetV2 and ResNet18-based visual encoders, respectively. Moreover, the BlazeNet64 visual encoder requires considerably less computation, with only 2.1 billion MACs per second of processing, in contrast to the 7.3 billion and 12.9 billion MACs required by the ShuffleNetV2 and ResNet18-based encoders, respectively.

In Table 2, we compare the performance of our visual encoder with the baselines on the simulated LRS3 dataset. In systems 1 to 3, where AV-ConvTasNet [4,31] is used as the speaker extractor, the ResNet18 visual encoder achieves the best overall SI-SNRi of 9.2 dB. Our visual encoder closely matches the ResNet18 visual encoder with an SI-SNRi of 9.1 dB, and surpasses the ShuffleNetV2 visual encoder by 0.3 dB. The AV-ConvTasNet speaker extractor shows suboptimal performance in the switching target scenario, with a notable drop in SI-SNRi after the switch. For example, system 3 has an SI-SNRi ‘after’ switching of 8.5 dB, which is 0.7 dB lower than the SI-SNRi ‘before’ switching.

In systems 4 to 6, where AV-SkiM [21] is used as the speaker extractor, the ResNet18 visual encoder achieves an SI-SNRi of 9.0 dB, while our visual encoder achieves a similar SI-SNRi of 9.1 dB. Notably, the AV-SkiM speaker extractor maintains consistent performance ‘before’ and ‘after’ the target switch in the switching target scenario, with similar SI-SNRi values. Across systems 1 to 6, the SI-SNRi in scenarios ‘without’ a target switch is generally higher than in scenarios with a target switch, due to the longer audio clips providing more context for speaker extraction.

In systems 7 and 8, we evaluate our visual encoder against the ResNet18-based visual encoder under non-causal (offline) conditions. Both systems achieve an average SI-SNRi of 12.8 dB, which is significantly higher than the performance of causal models. In systems 9 and 10, we test the impact of reducing the visual frame rate to 12.5 and 5 frames per second (FPS), respectively. While this reduction leads to lower computational costs, the performance drops significantly, indicating that the reduced frame rates are not worthwhile given the substantial loss in performance.

4.2. Study on the acoustic encoder

In table 3, we evaluate the effectiveness of our proposed acoustic encoder. System 6 serves as the baseline, which does not include the acoustic encoder, whereas System 12 incorporates our proposed visual encoder along with the time-frequency-domain hybrid loss function. The results indicate that System 12 outper-

Table 2: Comparison of our visual encoder with baselines on the LRS3 dataset. We compare its effectiveness with different speaker (spk.) extractor backbones, causal or non-causal settings, and different visual frame rates (V. FPS). We first report the SI-SNRi in dB, SNRi in dB, PESQi, and STOIi for ‘all’ the clips, we then report the SI-SNRi in dB for utterance segments ‘before’ the target switch, ‘after’ the target switch, and utterances ‘without’ target switch. The acoustic encoder is not used in these systems.

Sys	Spk. Extractor	Causal	V. FPS	V. Encoder	SI-SNRi (all)	SNRi (all)	PESQi (all)	STOIi (all)	SI-SNRi (before)	SI-SNRi (after)	SI-SNRi (without)	Params	MACs
1	AV-ConvTasNet	✓	25	ResNet18	9.2	9.5	0.70	0.17	9.4	8.5	9.9	22.1	31.2
2				ShuffleNetV2	8.8	9.1	0.65	0.16	8.9	8.2	9.3	11.8	25.6
3				BlazeNet64 (Ours)	9.1	9.4	0.68	0.17	9.2	8.5	9.7	11.0	20.8
4			25	ResNet18	9.0	9.3	0.72	0.16	8.7	8.9	9.4	19.1	18.2
5				ShuffleNetV2	9.0	9.3	0.69	0.17	9.0	8.8	9.5	8.8	12.7
6				BlazeNet64 (Ours)	9.1	9.4	0.72	0.17	8.9	8.9	9.7	8.1	7.9
7	AV-SkiM	✗	25	ResNet18	12.8	13.1	1.08	0.21	13.1	12.1	13.4	32.9	23.7
8				BlazeNet64 (Ours)	12.8	13.0	1.09	0.21	13.1	12.1	13.4	21.8	13.0
9		✓	12.5	BlazeNet64 (Ours)	6.3	6.6	0.46	0.12	6.2	5.8	6.8	8.1	6.6
10			5		0.0	0.0	-0.30	0.00	0.0	-0.1	0.0	8.1	6.0

Table 3: The studies of our proposed acoustic encoder on the LRS3 datasets are detailed in the table. We present results for configurations with (✓) and without (✗) the acoustic encoder, as well as for systems trained with different loss functions. The SI-SNRi and SNRi values are reported in dB. All systems evaluated in this table utilize our proposed BlazeNet64 visual encoder. Bolded values highlight the best results obtained.

Sys	Spk. Extractor	Acoustic Encoder	Loss	SI-SNRi (all)	SNRi (all)	PESQi (all)	STOIi (all)	SI-SNRi (before)	SI-SNRi (after)	SI-SNRi (without)	Params	MACs
6	AV-SkiM	✗	\mathcal{L}_{SI-SNR}	9.1	9.4	0.72	0.17	8.9	8.9	9.7	8.1	7.9
11	AV-SkiM-A	✓	\mathcal{L}_{SNR}	9.1	9.4	0.72	0.17	9.0	8.9	9.4	8.6	8.9
12			$\mathcal{L}_1 + \mathcal{L}_2$	10.0	10.4	0.83	0.18	9.9	9.9	10.5	8.6	8.9
13	AV-SkiM	✗	\mathcal{L}_1	9.5	9.7	0.75	0.18	9.3	9.5	9.7	8.1	7.9

Table 4: The SI-SNRi (dB) of models trained and evaluated on the VoxCeleb2 and TCD-TIMIT dataset mixtures.

Sys	Model	VoxCeleb2	TCD-TIMIT
14	SkiM w/ ResNet18	5.6	9.1
15	SkiM w/ ShuffleNet	5.5	8.5
16	SkiM w/ BlazeNet64 (Ours)	5.6	9.5
17	SkiM w/ BlazeNet64 & acoustic encoder (Ours)	6.4	11.0

forms System 6 by 0.9 dB in SI-SNRi, despite only a modest increase of 0.5 million parameters and 1 billion MACs. Notably, while the acoustic encoder is designed to leverage extraction momentum from past frames, it maintains performance stability even in scenarios where the target speaker changes. Specifically, the SI-SNRi values before and after a target change remain consistent at 9.9 dB, demonstrating the robustness of the network in adapting to changes and effectively resetting the extraction momentum.

We employed the PARIS [28] training strategy for our autoregressive acoustic encoder. However, instead of the SNR loss function used in the original PARIS paper, we utilized the time-frequency-domain hybrid loss function ($\mathcal{L}_1 + \mathcal{L}_2$) for training. In our evaluation, System 11 represents a baseline trained with the SNR loss function alongside the acoustic encoder. It is observed that System 11 performs comparably to System 6, which is trained using SI-SNR, and does not surpass our proposed System 12. This result aligns with the general finding that systems trained with the SNR loss function typically do not achieve the same level of performance as those trained with the SI-SNR loss function. Conversely, our proposed System 12, which integrates

SI-SNR loss with a frequency domain loss to regulate the output signal energy, demonstrates superior effectiveness.

We also present system 13 as an ablation study, which does not employ the acoustic encoder but is trained using the time-frequency-domain hybrid loss function \mathcal{L}_1 , it has SI-SNRi of 9.5 dB, which is better than system 6 but not better than system 12, showing the effectiveness of our acoustic encoder.

4.3. Comparison on more datasets

In Table 4, we compare our model and baselines on the VoxCeleb2 [41] and TCD-TIMIT dataset [42], where the simulated mixtures are from [31]. On both dataset, our system 16 performs comparably or outperforms baseline systems 14 and 15 in terms of SI-SNRi. Our system 17 with an acoustic encoder performs the best with SI-SNRi of 6.4 dB and 11.0 dB on VoxCeleb2 and TCD-TIMIT mixtures, respectively.

5. Conclusion

In conclusion, this work presents a significant advancement in online audio-visual speaker extraction by addressing both computational efficiency and performance. The proposed visual encoder, with its lightweight design and efficient processing, provides a competitive alternative to the more complex visual encoder. The novel acoustic encoder, leveraging past frame information, effectively enhances the extraction process, and is robust to the dynamic scenarios involving target speaker changes. Overall, our approach demonstrates strong performance across various metrics with practical advantages in terms of computational efficiency for real-time applications in multi-talker environments.

6. References

- [1] E. C. Cherry, “Some experiments on the recognition of speech, with one and with two ears,” *J. Acoust. Soc. Am.*, vol. 25, no. 5, pp. 975–979, 1953.
- [2] J. Wang, Z. Pan, M. Zhang, R. T. Tan, and H. Li, “Restoring speaking lips from occlusion for audio-visual speech recognition,” in *Proc. AAAI*, vol. 38, 2024.
- [3] J. R. Hershey, Z. Chen, J. Le Roux, and S. Watanabe, “Deep clustering: Discriminative embeddings for segmentation and separation,” in *Proc. ICASSP*, 2016, pp. 31–35.
- [4] Y. Luo and N. Mesgarani, “Conv-TasNet: Surpassing ideal time–frequency magnitude masking for speech separation,” *IEEE/ACM Trans. Audio, Speech, Lang. Process.*, vol. 27, no. 8, pp. 1256–1266, 2019.
- [5] Y. Luo, Z. Chen, and T. Yoshioka, “Dual-Path RNN: Efficient long sequence modeling for time-domain single-channel speech separation,” in *Proc. ICASSP*, 2020, pp. 46–50.
- [6] Z.-Q. Wang, S. Cornell, S. Choi, Y. Lee, B.-Y. Kim, and S. Watanabe, “TF-GridNet: Making time-frequency domain models great again for monaural speaker separation,” in *Proc. ICASSP*, 2023.
- [7] S. Zhao, Y. Ma, C. Ni, C. Zhang, H. Wang, T. H. Nguyen, K. Zhou, J. Yip, D. Ng, and B. Ma, “Mossformer2: Combining transformer and rnn-free recurrent network for enhanced time-domain monaural speech separation,” in *Proc. ICASSP*, 2024.
- [8] Q. Wang, H. Muckenhirn, K. Wilson, P. Sridhar, Z. Wu, J. R. Hershey, R. A. Saurous, R. J. Weiss, Y. Jia, and I. L. Moreno, “Voice-Filter: Targeted voice separation by speaker-conditioned spectrogram masking,” in *Proc. Interspeech*, 2019, pp. 2728–2732.
- [9] K. Žmolíková, M. Delcroix, K. Kinoshita, T. Ochiai, T. Nakatani, L. Burget, and J. Černocký, “SpeakerBeam: Speaker aware neural network for target speaker extraction in speech mixtures,” *IEEE J. Sel. Top. Signal Process.*, vol. 13, no. 4, pp. 800–814, 2019.
- [10] C. Xu, W. Rao, E. S. Chng, and H. Li, “SpEx: Multi-scale time domain speaker extraction network,” *IEEE/ACM Trans. Audio, Speech, Lang. Process.*, vol. 28, pp. 1370–1384, 2020.
- [11] A. Ephrat, I. Mosseri, O. Lang, T. Dekel, K. Wilson, A. Hassidim, W. T. Freeman, and M. Rubinstein, “Looking to listen at the cocktail party: a speaker-independent audio-visual model for speech separation,” *ACM Trans. Graph.*, vol. 37, no. 4, pp. 1–11, 2018.
- [12] Z. Pan, G. Wichern, Y. Masuyama, F. G. Germain, S. Khurana, C. Hori, and J. Le Roux, “Scenario-aware audio-visual TF-Gridnet for target speech extraction,” in *Proc. ASRU*, 2023.
- [13] J. Wu, Y. Xu, S. Zhang, L. Chen, M. Yu, L. Xie, and D. Yu, “Time domain audio visual speech separation,” in *Proc. ASRU*, 2019, pp. 667–673.
- [14] Z. Pan, R. Tao, C. Xu, and H. Li, “Selective listening by synchronizing speech with lips,” *IEEE/ACM Trans. Audio, Speech, Lang. Process.*, vol. 30, pp. 1650–1664, 2022.
- [15] S.-W. Chung, S. Choe, J. S. Chung, and H.-G. Kang, “FaceFilter: Audio-visual speech separation using still images,” in *Proc. Interspeech*, 2020, pp. 3481–3485.
- [16] E. Ceolini, J. Hjortkjær, D. D. Wong, J. O’Sullivan, V. S. Raghavan, J. Herrero, A. D. Mehta, S.-C. Liu, and N. Mesgarani, “Brain-informed speech separation (BISS) for enhancement of target speaker in multitalker speech perception,” *NeuroImage*, vol. 223, p. 117282, 2020.
- [17] Z. Pan, M. Borsdorf, S. Cai, T. Schultz, and H. Li, “NeuroHeed: Neuro-steered speaker extraction using EEG signals,” *arXiv preprint arXiv:2307.14303*, 2023.
- [18] Z. Pan, G. Wichern, F. G. Germain, S. Khurana, and J. Le Roux, “NeuroHeed+: Improving neuro-steered speaker extraction with joint auditory attention detection,” in *Proc. ICASSP*, 2024, pp. 11 456–11 460.
- [19] Q. Wang, I. L. Moreno, M. Saglam, K. Wilson, A. Chiao, R. Liu, Y. He, W. Li, J. Pelecanos, M. Nika *et al.*, “VoiceFilter-Lite: Streaming targeted voice separation for on-device speech recognition,” *Proc. Interspeech*, pp. 2677–2681, 2020.
- [20] H. Oh, J. Yi, and Y. Lee, “Papez: Resource-efficient speech separation with auditory working memory,” in *Proc. ICASSP*, 2023.
- [21] C. Li, L. Yang, W. Wang, and Y. Qian, “SkIM: Skipping memory LSTM for low-latency real-time continuous speech separation,” in *Proc. ICASSP*, 2022, pp. 681–685.
- [22] L. D. Libera, C. Subakan, M. Ravanelli, S. Cornell, F. Lepoutre, and F. Grondin, “Resource-efficient separation transformer,” in *Proc. ICASSP*, 2024.
- [23] H. Chen, R. Mira, S. Petridis, and M. Pantic, “RT-LA-VocE: Real-time low-SNR audio-visual speech enhancement,” *Proc. Interspeech*, 2024.
- [24] Z. Pan, M. Ge, and H. Li, “USEV: Universal speaker extraction with visual cue,” *IEEE/ACM Trans. Audio, Speech, Lang. Process.*, vol. 30, pp. 3032–3045, 2022.
- [25] Z. Zhu, H. Yang, M. Tang, Z. Yang, S. E. Eskimez, and H. Wang, “Real-time audio-visual end-to-end speech enhancement,” in *Proc. ICASSP*, 2023.
- [26] V. Bazarevsky, Y. Kartynnik, A. Vakunov, K. Raveendran, and M. Grundmann, “BlazeFace: Sub-millisecond neural face detection on mobile GPUs,” in *CVPR Workshop on Computer Vision for Augmented and Virtual Reality*, 2019.
- [27] F. Chollet, “Xception: Deep learning with depthwise separable convolutions,” in *Proc. CVPR*, 2017, pp. 1800–1807.
- [28] Z. Pan, G. Wichern, F. G. Germain, K. Saijo, and J. Le Roux, “PARIS: Pseudo-autoregressive siamese training for online speech separation,” in *Proc. Interspeech*, 2024.
- [29] T. Afouras, J. S. Chung, and A. Zisserman, “LRS3-TED: a large-scale dataset for visual speech recognition,” *preprint arXiv:1809.00496*, 2018.
- [30] —, “Deep lip reading: A comparison of models and an online application,” in *Proc. Interspeech*, 2018, pp. 3514–3518.
- [31] Z. Pan, R. Tao, C. Xu, and H. Li, “MuSE: Multi-modal target speaker extraction with visual cues,” in *Proc. ICASSP*, 2021, pp. 6678–6682.
- [32] J. Le Roux, S. Wisdom, H. Erdogan, and J. R. Hershey, “SDR–half-baked or well done?” in *Proc. ICASSP*, 2019, pp. 626–630.
- [33] Z. Pan, M. Ge, and H. Li, “A hybrid continuity loss to reduce over-suppression for time-domain target speaker extraction,” in *Proc. Interspeech*, 2022, pp. 1786–1790.
- [34] W.-N. Hsu, T. Remez, B. Shi, J. Donley, and Y. Adi, “ReVISE: Self-supervised speech resynthesis with visual input for universal and generalized speech regeneration,” in *Proc. CVPR*, 06 2023, pp. 18 795–18 805.
- [35] Q. Liu, M. Ge, Z. Wu, and H. Li, “PIAVE: A pose-invariant audio-visual speaker extraction network,” in *Proc. Interspeech*, 2023.
- [36] H. Martel, J. Richter, K. Li, X. Hu, and T. Gerkmann, “Audio-visual speech separation in noisy environments with a lightweight iterative model,” in *Proc. Interspeech*, 2023.
- [37] D. P. Kingma and J. Ba, “Adam, a method for stochastic optimization,” in *Proc. ICLR*, vol. 1412, 2015.
- [38] E. Vincent, R. Gribonval, and C. Févotte, “Performance measurement in blind audio source separation,” *IEEE Trans. Audio, Speech, Lang. Process.*, vol. 14, no. 4, pp. 1462–1469, 2006.
- [39] A. W. Rix, J. G. Beerends, M. P. Hollier, and A. P. Hekstra, “Perceptual evaluation of speech quality (PESQ)-a new method for speech quality assessment of telephone networks and codecs,” in *Proc. ICASSP*, 2001, pp. 749–752.
- [40] C. H. Taal, R. C. Hendriks, R. Heusdens, and J. Jensen, “A short-time objective intelligibility measure for time-frequency weighted noisy speech,” in *Proc. ICASSP*, 2010, pp. 4214–4217.
- [41] J. S. Chung, A. Nagrani, and A. Zisserman, “VoxCeleb2: Deep speaker recognition,” *Proc. Interspeech*, pp. 1086–1090, 2018.
- [42] N. Harte and E. Gillen, “TCD-TIMIT: An audio-visual corpus of continuous speech,” *IEEE Transactions on Multimedia*, vol. 17, no. 5, pp. 603–615, 2015.